

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-169172

(43) 公開日 平成11年(1999) 6月29日

(51) IntCl.⁶

識別記号

F I

C 1 2 N 15/09

C 1 2 N 15/00

A

G 0 1 N 33/50

G 0 1 N 33/50

P

審査請求 未請求 請求項の数11 O L (全 10 頁)

(21) 出願番号 特願平9-336858

(22) 出願日 平成9年(1997)12月8日

(71) 出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72) 発明者 村上 勝彦

東京都国分寺市東恋ヶ窪一丁目280番地

株式会社日立製作所中央研究所内

(72) 発明者 高木 利久

東京都港区白金台4-6-1

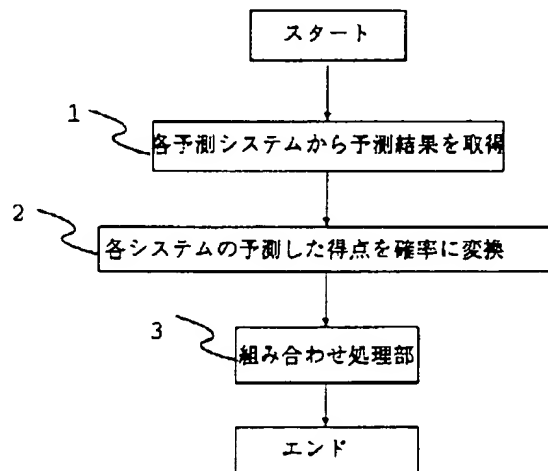
(74) 代理人 弁理士 平木 祐輔

(54) 【発明の名称】 DNA塩基配列上のタンパク質コード領域予測方法及び記録媒体

(57) 【要約】

【課題】 本発明の課題は、DNA塩基配列上のタンパク質コード領域を正確に予測するコード領域予測方法を提供すること。

【解決手段】 異なるアルゴリズムを用いた複数のコード領域予測方法の結果を入力として1、あらためてコード領域を予測するコード領域予測方法であって、各予測方法から出力されるスコアを予測領域が正解である確率に変換し2、確率の値を比較する3ことで信頼度の高い領域を選択して予測する。これによって、各方法に用いられているアルゴリズムを組み合わせることが容易にでき、その結果正解率を上げることができる。



【特許請求の範囲】

【請求項1】 DNA塩基配列データ上でタンパク質をコードする領域（これをコード領域と呼ぶ）を検出する方法において、異なるアルゴリズムを用いた複数のコード領域予測方法の予測結果を入力として、あらかじめコード領域を予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項2】 請求項1に記載のDNA塩基配列上のタンパク質コード領域予測方法において、各予測方法が共通にコード領域と予測した領域をとり、その領域をコード領域と予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項3】 請求項1に記載のDNA塩基配列上のタンパク質コード領域予測方法において、各予測方法のいずれかが予測した領域のすべてをとり、その領域をコード領域と予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項4】 請求項1に記載のDNA塩基配列上のタンパク質コード領域予測方法において、あらかじめ各予測方法に優先順位をつけておき、もし、重なる領域を複数の予測方法がコード領域と予測した場合、その境界の決定については、優先順位の高い予測方法の予測による境界を採用してコード領域と予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項5】 請求項1乃至4のいずれかに記載のDNA塩基配列上のタンパク質コード領域予測方法において、各コード領域予測方法がコード領域と予測した領域に与えられたスコアを、あらかじめ定めた関数によってその領域が正しくコード領域である確率に変換し、その確率の平均値があらかじめ定めた閾値よりも大きい時にその領域をコード領域と予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項6】 請求項1に記載のDNA塩基配列上のタンパク質コード領域予測方法において、もし、重なる領域を複数の予測方法がコード領域と予測した場合、その境界の決定については各コード領域予測方法がコード領域と予測した領域に与えられたスコアをあらかじめ定めた関数によってその領域が正しくコード領域である確率に変換し、確率が高い予測方法の予測による境界を採用してコード領域と予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項7】 請求項6に記載のDNA塩基配列上のタンパク質コード領域予測方法において、その選択したコード領域予測方法の確率があらかじめ定めた閾値よりも大きい時にその領域をコード領域と予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項8】 請求項6又は7に記載のDNA塩基配列上のタンパク質コード領域予測方法において、各コード領域予測方法がコード領域と予測した領域に与えたスコ

アを確率に変換する関数が、コード領域予測方法と、その方法が与えるスコアと、予測した領域の5'側及び3'側の境界の種類の関数であることを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項9】 請求項5、7又は8に記載のDNA塩基配列上のタンパク質コード領域予測方法において、各予測方法の確率の平均値をその領域の最終スコアとするものであって、コード領域という予測を陽性、非コード領域という予測を陰性として、TPを真陽性の塩基数、TNを真陰性の塩基数、FPを偽陽性の塩基数、FNを偽陰性の塩基数、PPを陽性の塩基数、PNを陰性の塩基数、APをコード領域全体の塩基数、ANを非コード領域全体の塩基数として、 $(TP)(TN) - (FP)(FN)$ を $(PP)(PN)(AP)(AN)$ の平方根で割ったものとして定義されるファイ相関係数を計算した時に、この相関係数が既知の配列データを解析した場合に最大になるようなスコアの閾値をもうけておき、領域の最終スコアがこの閾値以上の場合に、その領域をコード領域と予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項10】 請求項5、7又は8に記載のDNA塩基配列上のタンパク質コード領域予測方法において、各予測方法の確率の平均値をその領域の最終スコアとして、FPを偽陽性の塩基数、FNを偽陰性の塩基数、APをコード領域の塩基数、ANを非コード領域の塩基数とすると、 $E = (FN/AP + FP/AN) / 2$ によって定義される平均誤り率が既知のデータを解析した場合に最低になるような閾値をもうけ、領域の最終スコアがこの閾値以上をコード領域と予測することを特徴とするDNA塩基配列上のタンパク質コード領域予測方法。

【請求項11】 請求項1乃至10のいずれかに記載されたDNA塩基配列上のタンパク質コード領域予測方法をコンピュータに実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明はDNA塩基配列の情報処理、特に配列上のタンパク質コード領域を検出する方法に関する。

【0002】

【従来の技術】現在、ゲノム解析計画によってDNA配列が大量に決定されている。何の情報も付加されていないDNA塩基配列のデータに対して、その中でタンパク質がコードされている領域（コード領域）を予測することは、新しい遺伝子を発見し、医学、薬学的な研究をすすめる上で重要である。配列を決定した後で情報処理により遺伝子領域を推定できれば効率的に新しい遺伝子特定することができる。DNA配列中のコード領域を発見する方法としては、既知の核酸配列データベースに対して似た配列があるかどうかを検索する方法がある。新

しい配列の中で、既知の遺伝子に似ている配列は遺伝子である可能性が高いからである。しかし、長い領域にわたって問題のDNA配列と似ている遺伝子がデータベースにない場合にはこの方法は有効でない。近年、GRAIL (Proceedings, The Second International Conference on Intelligent Systems for Molecular Biology, page 376-384, 1994, ISBN 0-929280-68-7) という方法などのように、統計情報をもとにしてコード領域を見つける方法が進展してきた。

【0003】このGRAILで使われるコード領域予測方法のアルゴリズムは、主にコード領域に現れる数塩基の配列の統計、すなわち6塩基の短い配列(以下、6mer

$$CP_1 = \prod_{i=0}^{k-2} P(a_{3i+1} \dots a_{3i+5} a_{3i+6}) P(a_{3i+2} \dots a_{3i+6} a_{3i+7}) P(a_{3i+3} \dots a_{3i+7} a_{3i+8})$$

----- (数1)

この数値をコーディングポテンシャルと言い、その値の大きさはその領域がコード領域らしい程度を示す。これ以外にもいくつか類似の指標を計算し、それらをフィードフォワード型のニューラルネットワークに入力することにより、最終的にその領域に0以上1以下の得点を与え、一定値以上であればコード領域であると予測する。

【0005】他にも、同様にDNA配列からコード領域を予測する方法として、FEXH (Nucleic Acid Research, vol. 22, Num. 24, page 5156-5163, 1994)やGen

$$P(X) = \frac{F_c(X)}{F_c(X) + F_n(X)} \text{----- (数2)}$$

ここで、 $F_c(X)$ は、8merXがデータベース中のコード領域に出て来た頻度、 $F_n(X)$ は、8merXがデータベース中の非コード領域に出て来た頻度である。FEXHにおけるコーディングポテンシャルの計算は(数2)を用いて定義される。すなわち、任意の長さの配列について

$$CP_2 = \frac{1}{N-7} \sum_{i=1}^{N-7} P(a_i \dots a_{i+7}) \text{----- (数3)}$$

FEXHで使われるコード領域予測方法におけるコード領域の境界のシグナル検出方法について説明するが、その前に境界に関する基本的事柄を述べる。境界から5'側(左側)の塩基の位置を-1, -2, ...と表現し、境界からすぐ3'側(右側)の塩基の位置を境界に近いほうから1, 2, ...と表現する。0の位置はない。以下では例えば、-4から+3までの塩基位置の7文字からなる領域は(-4, 3)と表す。コード領域の境界では、境界からの位置によって使われる塩基の種類に偏りがあり、よく出て来る配列パターンを境界の「コンセンサス配列」という。コード領域の境界の種類はいくつかあって、原核生物の場合は開始コドンとよばれる配列'ATG'で始まるか、終止コドンと呼ばれる配列で終わるが、ヒトを初めとする真核生物の場合はさらに、コード領域の5'側の境界と3'側の境界とで2種類あり、この2つはそれぞれアクセプター

と呼ぶ)がコード領域に出現する頻度に基づいている。ある6mer (例えば、X=acgttc)がデータベース中の配列に出現した頻度のうち、コード領域に出て来た頻度 $F_c(X)$ と非コード領域に出て来た頻度 $F_n(X)$ の比 $P(X) = F_c(X)/F_n(X)$ をもって、この6merのスコアとする。さらに、6merでなく、ある程度長い領域がコード領域であるかどうかを識別するには、その領域に現れた6merすべてについて、上記の6merのスコアを(数1)にあるように乗じて、この領域のスコアとする。

【0004】

【数1】

eParser (Journal of Molecular Biology, 248, page 1-18, 1995)などがある。FEXHで使われるコード領域予測方法のアルゴリズムでは、GRAILとは異なるコーディングポテンシャルを用い、さらにコード領域の境界のシグナル検出方法を用いている。FEXHでは8merのスコアを(数2)によって計算する。

【0006】

【数2】

て、その配列のコーディングポテンシャルは、(数2)を考える領域にわたって平均したもの(数3)である。

【0007】

【数3】

サイト、ドナーサイトと呼ばれる。この種類によって、コンセンサス配列が違う。ほとんどの場合、コード領域の間に存在する介在配列と呼ばれる領域は、GTの2塩基で始まり、AGの2塩基で終わるので、配列GTが見つかればアクセプターサイトの候補であり、配列にAGが見つかればドナーサイトの候補である。これらのうち、実際にはコード領域の境界ではない位置を疑似境界部位と呼ぶ。

【0008】FEXHにおいては、コード領域のシグナルを検出するために、以下で定義される「3merのスコア」を用いる。まず、領域(L, R)に注目する。L, Rは領域の境界の位置である。ある3mer Y_k ($k=1, \dots, 64$)が領域(L, R)内の位置iにでてくる頻度を、実際の境界のデータと疑似部位に対してそれぞれ $F_{i,k}^L$, $F_{i,k}^R$ とおき、「位置iでの3merのスコア」を(数4)

によって定義する。

【0009】

【数4】

$$P(i) = \frac{F_{s,k}^i}{F_{s,k}^i + F_{p,k}^i} \dots\dots\dots (数4)$$

この(数4)を領域(L,R)に渡って平均したもの
(数5)が領域(L,R)の「3merのスコア」である。

【0010】

【数5】

$$P = \frac{1}{m} (\sum_{i=L}^R P(i)) \dots\dots\dots (数5)$$

ここで、mは領域(L,R)の長さである。さらに、
(数5)の和をとるときに、ある値 $\alpha=0.65$ を与えてお
き、 $P(i)$ が α よりも大きい $P(i)$ の和をとったときの
スコアを特に、「有意な3merのスコア」と定義する。さ
て、FEXHで使われるコード領域予測方法におけるコ
ード領域の境界のシグナル検出方法について説明する。
ドナーサイトについては、(-30,-5)のコーディングポ
テンシャル、領域(-4,6)の3merのスコア、領域(-30,-5)の
3merのスコア、領域(7,50)の3merのスコア、領域(-4,6)
の有意な3merのスコア、領域(6,50)のG, GG, GGG
の数などを計算し、統計的手法である判別分析によりこ
れらのスコアを組み合わせて、一つのドナーサイトに対
して一つの統合スコアを計算する。アクセプターサイト
については、領域(-48,-34)の3merのスコア、領域(-3
3,-7)の3merのスコア、領域(-6,5)の3merのスコア、領

域(6,30)の3merのスコア、領域(1,54)の8merのコー
ディングポテンシャル、領域(-1,-54)の8merのコー
ディングポテンシャル、領域(-33,-7)のT, Cの数を判別分析
によって組み合わせて、一つのアクセプターサイトに対
して一つの統合スコアを計算する。

【0011】ある領域がコード領域かどうかを決定す
ためのスコアは、コーディングポテンシャル、ドナーサ
イトのスコア、アクセプターサイトのスコアなどを判別
分析で組み合わせて、統合スコアが閾値より高い領域を
コード領域と予測する。GeneParser で使われるコ
ード領域予測方法のアルゴリズムでは、i 塩基目からj 塩
基目までの領域のコーディングポテンシャルを計算する
のに、(数6)を用いている。

【0012】

【数6】

$$IFhex(i, j) = \max \begin{cases} \sum_{k=0.3\dots j-6} \ln(\frac{f_k}{F_k}) \\ \sum_{k=1.4\dots j-9} \ln(\frac{f_k}{F_k}) \\ \sum_{k=2.5\dots j-12} \ln(\frac{f_k}{F_k}) \end{cases} \dots\dots\dots (数6)$$

ここで f_k は6mer $a_{i+k}, a_{i+k+1}, \dots, a_{i+k+5}$ ($a_1 \in \{A, C, G, T\}$) が学習データのコード領域に現れた頻度
で、6merの第一番目の文字がコドンの一番目になってい
るときだけ数えたものである。 F_k は、その配列と同じ
塩基組成でランダムに配列を生成した場合にその6merが

出て来る頻度の期待値である。

【0013】さらに、(数7)で定義される局所的複雑
度も用いている。

【0014】

【数7】

$$H = - \sum_{k \in \{A, C, G, T\}} (\frac{N_k}{L}) \log_2 (\frac{N_k}{L}) \dots\dots\dots (数7)$$

ここで、Lは考慮している配列の長さで、 N_k は塩基k
(kは、A, C, G, Tのいずれか) が長さLの配列に
現れた頻度である。GeneParser では、コード領域の
境界のスコアを以下のように計算する。境界の領域(i,

j) 内の配列 $s_i s_{i+1} \dots s_j$ に対してそのスコア $S(i, j)$
は(数8)で定義する。

【0015】

【数8】

$$S(i, j) = \sum_{k=i}^j \log_2 (f_{b,k} / p_b) \dots\dots\dots (数8)$$

ここで、 $f_{b,i}$ は、位置iに塩基b (bはA, C, G,
Tのいずれか) が出てきた頻度であり、 p_b は塩基bが
その配列に出てくる事前確率である。他にもコード領域
らしいかどうかのスコアや、ドナーサイトのスコア、ア
クセプターサイトのスコアを計算し、それらをフィード
フォワード型のニューラルネットワークで統合して、最

最終的にその領域に0以上1以下の得点を与え、一定値
以上であればコード領域と予測結果を出す。GeneParser
は、GRAILと同じくニューラルネットワークを用
いているが、先に述べたように考慮している特徴が異
なり、さらにその学習方法も異なっている。

【0016】このように、これらの予測方法は異なる特

微検出方法を用いており、それらのスコアを統合化する方法も異なっている。これらの方法の予測は不正解であることも多く、同じ配列を解析しても、検出しにくいコード領域に対しては方法によって予測結果が異なることが多い。

【0017】

【発明が解決しようとする課題】これまでは単一の方法をもとにしていたため、精度が低かった。これは、各予測方法の捉えている特徴が部分的なものだからである。したがって、できるだけ多くの特徴をとらえた結果をまとめて予測をして、正解率を上げることが課題である。一方、各方法で採用している配列の特徴検出方法を一つの方法の中で実装するのは、人的コストがかかる。そこで、出来るだけ簡便な方法で多くの特徴をとらえることが課題である。本発明の目的は、上記の課題を解決し、信頼性の高いコード領域予測をする方法を提供することが目的である。

【0018】

【課題を解決するための手段】上記課題を解決するために本発明では、DNA配列上にある遺伝子の様々な特徴を異なるアルゴリズムによって学習した遺伝子予測方法の結果を入力に用いて、それらを統合した予測を行う。これによって、一つの予測方法で考慮できない多くの特徴を考慮した予測結果を容易に得られる。また、異なる予測方法による予測結果の信頼性を比較したうえで最終予測を行うために、各予測方法のスコアを正解率に変換し、正解率を比較することによって、より信頼の高い予測結果を得る。

【0019】

【発明の実施の形態】図1は本発明の実施の概要で、複数の予測プログラムの結果を使って、新たな予測をする方法の流れ図である。各予測方法による予測処理部1で、入力されたDNA配列データを各コード領域予測方法で解析し、各々の予測結果を得る。各コード領域予測プログラムの出力結果の中から、入力配列の何塩基目から何塩基目までがコード領域であるかという情報と、そのスコア(確からしさ)の情報を保持しておく。次に、スコア変換部2で、各予測された領域のスコアからその領域が正解である確率を求める。正解であるとは、予測した領域と実際のコード領域がオーバーラップしていることである。この確率は、あらかじめ設定してある変換関数で得られる。この確率は、以下Pscoreと呼ぶ。変換関数の作成方法については、後述する。

【0020】各予測結果をもとにして、あらためてコード領域を予測する組み合わせ処理3を行い、最終予測結果を出力して終了である。図2はこの組み合わせ処理3の詳細についての説明図である。図3は複数の予測方法によって、異なる領域がコード領域と予測された例である。図3の上部に、3つのコード領域予測方法によってコード領域と予測された領域を描いた。すなわち、FE

XHによってコード領域と予測された領域21、GeneParserによってコード領域と予測された領域22、GRAILによってコード領域と予測された領域23である。それぞれのPscoreは、0.8、0.4、0.9である。横軸は入力したDNA配列上での塩基の位置を示す。まず、記録配列24を用意し、その全領域に整数0を入れて初期化11を行う。次に、各予測方法で予測された各塩基に対応する記録配列24の部分に1を足して予測結果を記録する処理12を行う。これによってどの場所がいくつの予測方法によってコード領域と予測されたかが、記録配列24に記録される。図3では、重なった予測方法の数を四角の高さで表しているが、実際には0以上の整数が入っている。

【0021】次に、スキャン処理13によって、記録配列24をスキャンし、1以上が記録されている重なり領域29をみつける。重なり領域29がなければ条件式14によって終了し、あればコード領域の境界を決定する処理15に進む。境界決定処理15と、最終スコアの決定16に関して、5つの方法を述べる。1つめの方法では、使用した全ての予測方法が予測した領域をコード領域の候補とする。このときの最終スコアは、各方法のPscoreの平均とする。この例では、最終スコアは $(0.8+0.4+0.9)/3 = 0.7$ である。以下これを「AND法」と呼ぶ。この方法によってコード領域と予測され得る領域25を図3に示した。この時点では、まだ候補であって、コード領域と予測されたわけではない。2つめの方法では、使用した予測方法のいずれかが予測した領域をコード領域の候補とする。このときの最終スコアは、各方法のスコアから換算した各Pscoreの平均とする。もし、いくつかの予測方法がその重なり領域中のどこもコード領域と予測していなかった場合の最終スコアは、予測しなかった予測方法のPscoreを0として、計算する。この例では、最終スコアは $(0.8+0.4+0.9)/3 = 0.7$ である。以下これを「OR法」と呼ぶ。この方法によってコード領域と予測され得る領域26を図3に示した。

【0022】3つめの方法では、その重なり領域29をコード領域と予測した予測方法のうち、最も高いPscoreを持っている予測方法を選択し、その予測方法が予測した領域をコード領域の候補とする。このときの最終スコアは、OR法と同様に各Pscoreの平均とする。もし、いくつかの予測方法がその重なり領域中のどこもコード領域と予測していなかった場合の最終スコアは、予測しなかった予測方法のPscoreを0として、計算する。この例では、最終スコアは $(0.8+0.4+0.9)/3 = 0.7$ である。以下これを「HIGHEST法」と呼ぶ。この方法によってコード領域と予測され得る領域27を図3に示した。

【0023】4つめの方法では、その重なり領域30をコード領域と予測した予測方法のうち、あらかじめつけておいた優先順位の高い予測方法の結果を選択する。その

予測方法が予測した領域をコード領域の候補とする。優先順位は、あらかじめ各予測方法のテストをしておきその成績の良い順にする。例えば、バーセッタらの行った遺伝子予測プログラムのテスト(Genomics, 34, page 353-367, 1996)において、境界も正確に予測したエキソン(コード領域の単位)の数の割合が高い順に設定すればよい。この順は、順位の高いほうからFEXH, GeneParser, GRAILである。このときの最終スコアは、選択した予測方法のPscoreとする。この例では、最終スコアはFEXHのPscoreになるので0.8である。以下これを「RULE法」と呼ぶ。この方法によってコード領域と予測され得る領域28を図3に示した。

【0024】5つめの方法では、その重なり領域30それぞれに対して、Pscoreと境界のタイプを考慮して新しいスコアEscoreを計算する。境界のタイプとは、開始コドン、ドナーサイト、アクセプターサイト、終止コドンのうちのいずれかである。この4つを以下では、i, d,

$$\text{Escore} = \text{Pe}(l, \text{ps}) \text{Pe}(r, \text{ps}) \dots\dots\dots (\text{数9})$$

この例でFEXH, GeneParser, GRAILのEscoreは、それぞれ0.72, 0.48, 0.54である。このEscoreが最も高い予測方法が予測した領域をコード領域の候補とする。このときの最終スコアは、選択した予測方法のEscoreとする。この例では、最終スコアは0.8*0.9=0.72である。以下これを「EDGE法」と呼ぶ。この方法によってコード領域と予測される領域29を図3に示した。

【0026】各方法のいずれかで候補の領域と最終スコアを決定した後、その候補をコード領域として予測するかどうかを閾値によって決定する。すなわち、最終スコアと閾値とを比較する部分17によって最終スコアが高いかどうかを判断し、高ければこれをコード領域と予測して、領域を出力する処理18を行う。ここで閾値の設定方

$$CC = \frac{(TP)(TN) - (FP)(FN)}{\sqrt{(PP)(PN)(AP)(AN)}} \dots\dots\dots (\text{数10})$$

ただし、コード領域という予測を陽性、非コード領域という予測を陰性として、TPは真陽性の塩基数、TNは真陰性の塩基数、FPは偽陽性の塩基数、FNは偽陰性の塩基数、PPは陽性の塩基数、PNは陰性の塩基数、APはコード領域の塩基数、ANは非コード領域の塩基数とする。-1 ≤ CC ≤ 1であり、完全に正解であればCC = 1であり、完全に正解と逆の予測をしたときはCC = -1となる。ランダムな予測に対しては、ファイ相関係数の期待値は0である。

【0028】閾値の決定に際しては、各方法毎に、多くの配列データにおいてあらかじめ0から1までのいろいろな閾値でファイ相関係数の平均を求めて、CCが最高値をとるときの閾値を、その方法の閾値と決める。閾値の決定に関しては、上記の方法の他に第1種の誤り率E₁ = FN/APと、第2種の誤り率E₂ = FP/ANの平均E_{AV} = (E₁ + E₂) / 2 が最低値をとるときの閾値

a, tと表記する。Escoreは以下のように計算する。まず、ある予測方法が予測したコード領域のうち、コード領域のPscoreがpsで、境界のタイプtypeにおける正解率をPe(type, ps)とする。これは、各予測方法毎に学習データから近似関数を作成しておく。この例では、FEXHの境界のタイプがa, tでPscoreが0.8であるから、FEXHのP(a, 0.8)とP(t, 0.8)を計算しておく。また、GeneParserの境界のタイプがa, dでPscoreが0.4であるからGeneParserのP(a, 0.4)とP(d, 0.4)を計算しておく。さらに、GRAILの境界のタイプがi, tでPscoreが0.9であるから、GRAILのP(i, 0.9)とP(t, 0.9)を計算しておく。次に、ある予測方法が予測した領域の左右の境界のタイプがl及びrならば(l, rはi, d, a, tのいずれかである)、このときのEscoreは、(数9)と定義する。

【0025】

【数9】

法について述べる。閾値を設定するときに高く設定すると、偽陽性の数は減るが感度が悪くなる。逆に、閾値を低く設定すると、感度が上がるが偽陽性の数が増えてしまう。そこで、閾値は何らかの指標が最適になるような適当な値に決めなければならない。ここでは、多くの配列データに対して、予測と正解の相関を示すファイ相関係数の全データに渡っての平均が、最高になるように定める。このファイ相関係数は、正解の分かっている配列、すなわち、コード領域の位置が分かっているDNA配列一つに対して一つの値が求められる。また、このファイ相関係数の定義は、(数10)である。

【0027】

【数10】

にする方法も考えられる。あるいは、上記の誤り率を計算する際に塩基数でなくコード領域の数で計算した誤り率の平均を最低値にするような閾値を採用する方法も考えられる。

【0029】ここで、各予測方法のスコアを予測した領域が正解である確率に変換する関数の作成方法を図4に沿って述べる。多くのデータを解析する必要があるので、はじめにさまざまな条件を満たすデータのみを集める処理31を行う。例えば、核酸配列データベースであるGenBankリリース100(1997年4月)の中から、項目'SOURCE'が'Homo sapiens'であり、一つ以上の'CDS'を含んでいるDNA配列のエントリーを集める。これらはコード領域の位置が分かっているDNA配列データである。また、イントロン領域を含むものについては、イントロン領域がGTで始まるかAGで終ることが条件であるので、この条件を満たさないデータは

捨てる。また、pseudo, putative, ORF, alternative, predict, fusionのうちいずれかの記述が項目'CDS'の中にあれば、それは実験的にコード領域とは確認されていないか、確実なコード領域がわかっていない可能性が強いため、これを除く。さらに、各コード領域予測方法が学習に用いたデータをここで使わないようにするため、1996年6月より前に登録されたデータを除く。

【0030】これらの処理31を経てデータセット32を作成する。データ一つについて各予測方法での解析33を実行した後、各予測方法のスコアとエラー率の関係を調べ

$$Pscore(score) = A + B \times score \dots\dots\dots (数11)$$

本方法によって正解率がどの程度変わるかを示す。DNA配列データセット32を本方法で解析し、配列一つごとに以下の正解率(数12)を計算し、これをデータの本数で割ったものである。なお、この正解率の計算方法はパーセットらによって提案された正解率であり(Genomic

$$AC = \frac{1}{2} \left[\frac{TP}{TP+FN} + \frac{TP}{TP+FP} + \frac{TN}{TN+FP} + \frac{TN}{TN+FN} \right] - 1 \dots\dots (数12)$$

各予測方法単独のときの数6の値はFEXHが0.64、GeneParserが0.63、GRAILが0.67であったのに対し、3方法の組合せを本発明のように行くと、AND法では0.53と低くなったが、OR法で0.76、HIGHEST法で0.74、RULE法で0.71、EDGE法で0.74と高くなり、後者4つの方法では正解率が上がった。

【0033】AND法は全体の正解率では、単独の結果より悪いが実際のコード領域のうち、検出できなかったコード領域の率をみると、単独では、FEXHが0.47、GeneParserが0.45、GRAILが0.27であるのに対し、AND法による3つのコード領域予測方法を組合せ

(表1) 他のアルゴリズムを用いたコード領域予測方法

| 名前 | 方法 | 開発元 |
|------------|-------------|------------------|
| ER | 線形判別分析 | 東大(日本) |
| GENSCAN | マルコフモデル | スタンフォード大学 |
| GenLang | 確定筋文法 | ペンシルベニア大学 |
| GenView | 線形結合 | イタリア国立研究会議(伊) |
| GeneID | パーセプトロン | ボストン大学 |
| GeneMark | マルコフモデル | ジョージア工科大学 |
| Genie | 隠れマルコフモデル | カリフォルニア大学 |
| MORGAN | 決定木 | ジョンズホプキンス大学 |
| MZEF | 2次判別分析 | コールドスプリングハーバー研究所 |
| OC1 | 決定木 | ジョンズホプキンス大学 |
| PROCRUSTES | スプライスアライメント | 南カリフォルニア大学 |
| Sorfind | ルールベース | ラビットハッチ・コーポレーション |
| VEIL | 隠れマルコフモデル | ジョンズホプキンス大学 |

【0036】

【発明の効果】従来、一つのアルゴリズムによる場合は

る処理34をする。実際に調べたFEXHのヒストグラムを図5に、GeneParserのヒストグラムを図6に、GRAILのヒストグラムを図7に示す。これらのヒストグラムから、各予測方法ごとに(数11)で仮定した変換関数を求めるため最小自乗法によるパラメータ推定35を行う。こうして出来た関数が、求める変換関数である。この変換関数は、一次式でなく二次式でもよい。

【0031】

【数11】

s, 34, page 353-367, 1996)、広く使われているものである。

【0032】

【数12】

ると、0.07と低くなり、93%のコード領域を検出することができた。AND方法は、感度と特定度のバランスでは、単独のときの正解率におよばないが、特別な場合、すなわち偽陽性が多くてもコード領域として可能性のある領域を出来るだけ多くリストアップしたい場合には有効である。

【0034】なお、上記以外にも他のアルゴリズムを用いたコード領域予測方法が知られているので、(表1)に例示する。これらの予測方法も任意に選択して用いることができる。

【0035】

【表1】

少ない指標にもとづいて予測していたため、効率良く多くのコード領域を得ることができなかった。本発明によ

って、複数のアルゴリズムで多角的に候補を選定できるので偽陽性があまり増えずに効率的に多くのコード領域を検出できる。その結果、全体的に精度が上がる。

【図面の簡単な説明】

【図1】本発明におけるコード領域予測手順

【図2】組合せ処理の説明図

【図3】重なり領域での境界決定方法の説明図

【図4】スコアを確率へ変換する関数を作成する方法の説明図

【図5】予測方法FEXHのスコアとエラー率の関係を示す図である。

【図6】予測方法GeneParserのスコアとエラー率の関係を示す図である。

【図7】予測方法GRAILのスコアとエラー率の関係を示す図である。

【符号の説明】

- 1…DNA配列を各予測方法で解析する処理部
- 2…予測方法のスコアを正解である確率に変換する部分
- 3…組み合わせ処理部
- 11…記録配列の初期化
- 12…各方法の予測領域を記録配列に記録する処理
- 13…記録配列のスキャン処理
- 14…重なり領域のチェック

15…境界の決定

16…最終スコアの決定

17…最終スコアと閾値の比較

18…コード領域と判定した領域の出力

21…予測方法FEXHが予測したコード領域

22…予測方法GeneParserが予測したコード領域

23…予測方法GRAILが予測したコード領域

24…DNA配列に対応する記録配列

25…「AND法」によってコード領域と予測される領域

26…「OR法」によってコード領域と予測される領域

27…「HIGHEST法」によってコード領域と予測される領域

28…「RULE法」によってコード領域と予測される領域

29…「EDGE法」によってコード領域と予測される領域

30…各予測方法による予測の重なり領域

31…核酸データベースから条件にあうデータを取得する処理

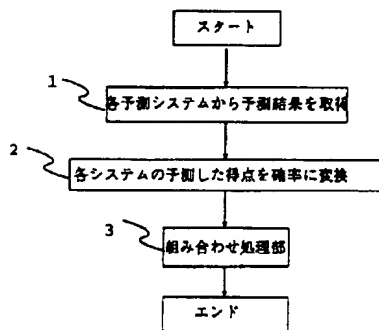
32…条件にあうデータの集合

33…各予測方法で解析する部分

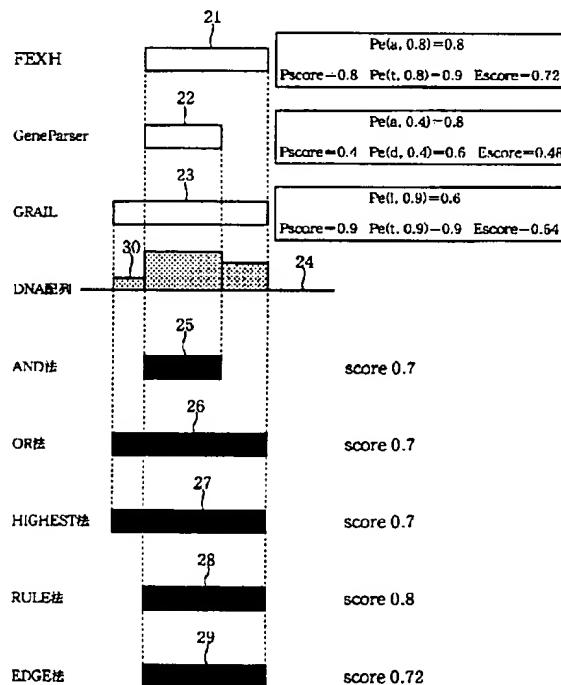
34…エラー率とスコアを計算する部分

35…変換関数のパラメータを計算する部分

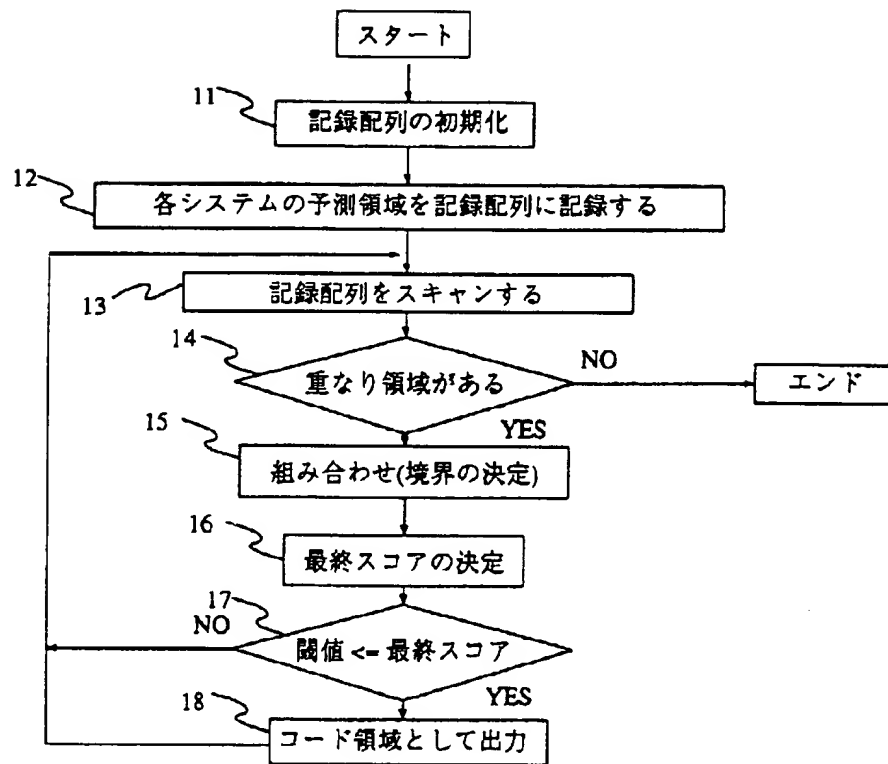
【図1】



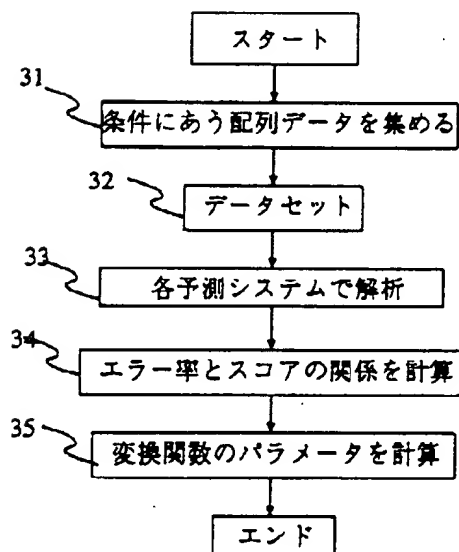
【図3】



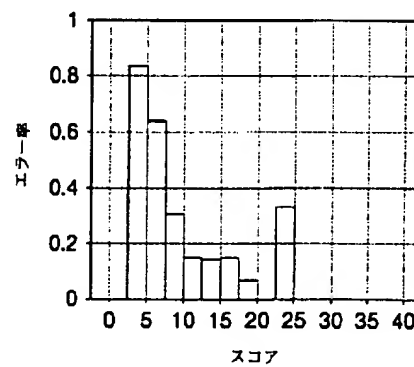
【図2】



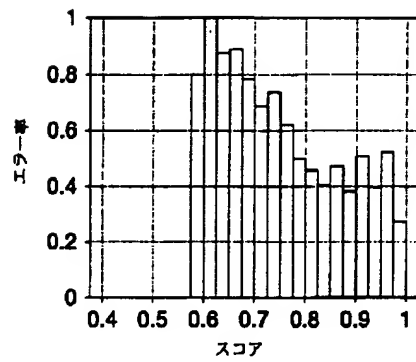
【図4】



【図5】



【図6】



【図7】

